

Bootstrap 方法及其应用

茅 宁

摘 要 1977年埃弗伦教授提出了一种新的统计推断方法——Bootstrap方法。该方法适用于一大类估计问题，因而受到越来越多的统计学者的重视。本文介绍了该方法的基本思想及其应用实例，并讨论了存在的问题及今后研究的方向。

一、引 言

1977年，美国Stanford大学统计系教授 Efron 在总结、归纳前人研究成果的基础上提出了一种新的统计方法——Bootstrap方法。最近几年来，该方法引起了学术界的重视，我国也已初步开展这方面的研究。从最初的定义看，Bootstrap方法是一种非参数统计方法，但实际上不仅仅局限于非参数的范畴，对于参数统计推断的某些问题同样也是适用的。从另一方面看，Bootstrap方法实质上是一种再抽样(resampling)过程。再抽样的方法还有许多，如Jackknife方法、无穷小 Jackknife方法、Delta方法、Half-Sampling方法等。本报告主要介绍 Bootstrap方法的基本思想及其若干应用问题，其中包括作者个人的一些看法，在理论上不准备做详细的论述。其它的再抽样方法可阅读报告中列出的有关文献。

下面先简单介绍一下问题的基本提法。设随机子样 $X = (X_1, \dots, X_n)$ 来自未知的总体分布 F ； $R(X, F)$ 为某个预先选定的随机变量，它是 X 和 F 的函数。现要求根据子样观测值 $X = x = (x_1, \dots, x_n)$ 估计 $R(X, F)$ 的分布特性，如均值、方差或分布密度函数等。例如，设 $\theta = \theta(F)$ 是总体分布 F 的某分布参数（如均值、方差、相关系数），而 $\hat{\theta}(X)$ 是关于 θ 的估计（如子样均值、子样方差、子样相关系数）。定义

$$R(X, F) = \hat{\theta}(X) - \theta(F) \quad (1.1)$$

现要由观测 $X = x$ 估计 $R(X, F)$ 的分布特性。注意此时 $R(X, F)$ 的均值和方差分别为 $\hat{\theta}$ 的偏和误差方差，正是我们所希望知道的。总体分布 F 一般是未知的，则上述问题为非参数问题。此外，也可能已知 F 的形式，但其某些分布参数未知，则此时是参数问题。

其次分析一下 Bootstrap方法产生的历史背景。从原则上讲，人们熟悉的参数统计推断方法是在计算能力受到严重束缚（计算费用昂贵、计算速度缓慢）这一特定的历史条件下产生和发展的。在处理问题时总认为总体分布形式已知，特别通常假定是正态分布。在此基础上建立起一套“墨守陈规”的模式，即对任何一个具体问题，总是利用现

成的理论和统计表(如 F 分布表, t 分布表)进行推断。实际上,对总体分布做出假定一般出于计算上的考虑,真实分布通常是无法精确了解的。有人曾经这样说:对于“为什么要采用正态假设”这一问题的最诚实回答是“只有在正态假设下才能获得问题的解”。事实也确实如此。可见用参数方法解决问题在相当程度上受到主观因素的影响。因此,从五十年代开始,非参数统计方法逐渐得到发展并受到重视。特别在电子计算机高度发展的今天,有必要也完全有可能从正态假设的约束中解脱出来,借助于计算机这一强有力的工具,创造出新的具有更高的灵活性、现实性和可靠性的统计推断方法。在解决具体问题时,希望不依赖于统计表,而是一切从零出发,只依据观测信息作分析和判断。Bootstrap 方法正是计算机时代的产物。Bootstrap 可译作“自助 (self-help)”,它说明该方法只依赖于给定的观测信息,不需要其它的假设和增加新的观测。

二、Bootstrap 方法的主要内容

1. Bootstrap 方法的基本步骤

首先重复一下问题的提法。考虑单子样场合,设随机子样 $X = (X_1, \dots, X_n)$ 来自未知的总体 F , $R(X, F)$ 为一预先选定的随机变量,它是 X 和 F 的函数。现要求根据观测 $x = (x_1, \dots, x_n)$ 去估计 $R(X, F)$ 的统计特性。

Bootstrap 方法实质上为一再抽样过程,其基本步骤如下:

(1) 由子样观测值构造子样经验分布函数 \hat{F}_n : \hat{F}_n 在每点 x_i 处具有质量 $\frac{1}{n}$, $i = 1, \dots, n$ 。 (2.1)

(2) 从 \hat{F}_n 中抽取简单子样 $X_i^* = x_i^*$, $i = 1, \dots, n$ 。称 $X^* = (X_1^*, \dots, X_n^*)$ 为 Bootstrap 子样, $X_i^* \sim \hat{F}_n$, $i = 1, \dots, n$ 。

(3) 用 $R^* = R(X^*, \hat{F}_n)$ 的分布来逼近 $R(X, F)$ 的分布。 R^* 的分布称为 Bootstrap 分布。

例如由 (1.1) 定义的 $R(X, F)$, 知此时

$$R^* = R(X^*, \hat{F}_n) = \hat{\theta}(X^*) - \theta(\hat{F}_n) \quad (2.2)$$

可见有两处做了近似处理:一是用 X^* 代替了 X , 二是用 \hat{F}_n 的分布参数代替了 F 的分布参数。

2. 几点说明

(1) 从非参数统计的观点看, \hat{F}_n 是 F 的非参数极大似然估计。它为一离散型分布,其可能值为 $x_i (i = 1, \dots, n)$, 且其均值和方差分别是

$$E_*[X^*] = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}, \quad (2.3)$$

$$Var_*[X^*] = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = s^2 \quad (2.4)$$

(2) 以上步骤是对一般非参数情况提出来的。若已知 F 的分布形式,则可对步骤 (1) 作适当修正,相应地称为参数 Bootstrap 方法。如若已知 F 为正态分布 $N(\mu, \sigma^2)$, 其中 μ 、 σ^2 未知,则可用 \bar{x} 和 s^2 代替未知参数 μ 和 σ^2 , 用 $N(\bar{x}, s^2)$ 代替 \hat{F}_n , 从中抽取 Boots-

trap子样。

(3) 若已知 F 为一连续分布, 为了适当考虑 F 的连续性, 可用 \hat{F}_n 与某一已知的连续分布 G 的卷积来代替 \hat{F}_n , 从中抽取 Bootstrap 子样, 相应地称为平滑 Bootstrap 方法, 即作

$$\hat{F}^c = \hat{F}_n * (cG) \quad (2.5)$$

其中 $c \in [0, 1]$ 。由 \hat{F}^c 中抽取 x_i^c , ($i=1, \dots, n$), 可这样得到: 首先分别由 \hat{F}_n 和 G 中抽取子样 x_i^* 和 y_i , ($i=1, \dots, n$), 再作

$$x_i^c = \sqrt{1-c^2} x_i^* + c y_i, \quad i=1, \dots, n \quad (2.6)$$

即可。特别有 $\hat{F}^0 = \hat{F}_n$, $\hat{F}^1 = G$ 。

(4) 设所考虑的随机变量 $R(X, F)$ 关于 X 对称, 即 X 的各分量 X 的排列不影响 $R(X, F)$ 的取值。定义再抽样矢量

$$P^* = (P_1^*, \dots, P_n^*) \quad (2.7)$$

它满足约束: $P_i^* \geq 0$, $\sum_{i=1}^n P_i^* = 1$ 。Bootstrap 方法即是考虑所有形如 $\frac{N^*}{n}$ 的 P^* , 其中

$$N^* = (N_1^*, \dots, N_n^*) \quad (2.8)$$

N_i^* 为 Bootstrap 子样 X^* 中等于 x_i 的分量 x_i^* 的个数, ($i=1, \dots, n$), $\sum_{i=1}^n N_i^* = n$ 是显然的。

由定义, $P^* = \frac{N^*}{n}$ 服从多项分布, 其均值和方差分别是

$$E[P^*] = \frac{e}{n} \quad (2.9)$$

$$V_{ar}[P^*] = \frac{1}{n^2} - \frac{e^T e}{n^3} \quad (2.10)$$

其中 $e = (1, 1, \dots, 1)$ 。由于在 $R(X, F)$ 为对称时, P^* 对于计算 $R(X^*, \hat{F}_n)$ 已足够了, 故可记

$$R^* = R(P^*) = R(X^*, \hat{F}_n) \quad (2.11)$$

(5) 在再抽样时, Bootstrap 子样 X^* 的容量可以不为 n , 一般有

$$X^* = (X_1^*, \dots, X_m^*), \quad X_i^* \sim \hat{F}_n \quad (2.12)$$

对于 $m \neq n$ 的场合有时具有重要的统计意义。例如可设 $m = n^2$ 或 $m = \sqrt{n}$, 这样可以研究分析判断的结论在子样容量变化时的情况。在下面的讨论中我们均考虑 $m = n$, 但这不是唯一选择。

3. 获得 Bootstrap 分布的途径

显然, 一旦观测子样 x 给定, \hat{F}_n 便可唯一确定, Bootstrap 子样也可得到, 剩下的问题在于如何得到 Bootstrap 分布。该分布实质上是在给定 x 之下的条件分布。根据不同情况, 可由以下途径获得。

途径 1 直接的理论计算

例 1 设 F 为 0-1 分布, 分布参数 $\theta(F) = P\{X=1\}$, 则

$$R(X, F) = \bar{X} - \theta(F), \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$$

给定观测 x 后, Bootstrap子样 $X^* = (X_1^*, \dots, X_n^*)$ 的每个分量 X_i^* 显然也服从 0-1 分布, 即

$$\begin{aligned} P_*\{X_i^* = 1\} &= \bar{x} = \theta(\hat{F}_n) \\ P_*\{X_i^* = 0\} &= 1 - \bar{x}, \quad i = 1, \dots, n \end{aligned}$$

“ $*$ ”表示概率运算关于 \hat{F}_n 进行。此时

$$R^* = R(X^*, \hat{F}_n) = X^* - \bar{x}, \quad \bar{X}^* = \frac{1}{n} \sum_{i=1}^n X_i^*$$

于是有

$$\begin{aligned} E_*[R^*] &= E_*[\bar{X}^* - \bar{x}] = \bar{x} - \bar{x} = 0 \\ \text{Var}_*[R^*] &= E_*[(\bar{X}^* - \bar{x})^2] = \frac{1}{n} \bar{x}(1 - \bar{x}) \end{aligned}$$

\bar{x} 作为观测的函数已给定。该结果与 Bootstrap 子样无关。

途径 2 用 Monte—Carlo 方法做逼近

在大多数场合直接计算是不可能的, 因此必须借助于计算机这一有力工具进行统计模拟计算。即从 \hat{F}_n 中重复抽取 N 个 Bootstrap 子样 $X^*(1), \dots, X^*(N)$; 计算相应的 $R^*(1) = R(X^*(1), \hat{F}_n), \dots, R^*(N) = R(X^*(N), \hat{F}_n)$; 从而用 $R^*(i), i = 1, \dots, N$ 的频率曲线 (即直方图) 作为对 Bootstrap 分布的逼近, 或者对 $R(X^*, \hat{F}_n)$ 的统计特性 (一、二阶矩等) 进行估计。

例 2 考虑参数估计的均方误差。设 $f(x; \theta)$ 为总体的分布密度, 其形式已知。则由观测 x 可对未知参数 θ 进行最大似然估计 $\hat{\theta}_{ML}$ 。同时, 也希望知道该估计的均方误差

$$\sigma^2 = E[(\theta - \hat{\theta}_{ML})^2] \quad (2.13)$$

当然 σ^2 本身也可用 ML 方法估计, 但在大多数场合很难得到问题的解。另一个可能的途径是由 Fisher 信息限给出 σ^2 的下界, 但这通常也是十分困难的。用 Bootstrap 方法解决上述问题却十分简单:

(1) 用 $f(x^*; \hat{\theta}_{ML})$ 代替 $f(x; \theta)$, 从中抽取 N 个相互独立的简单子样 $x^*(i), i = 1, \dots, N$, 其中

$$x^*(i) = (x_1^*(i), \dots, x_n^*(i))$$

同时计算相应的 $\hat{\theta}_{ML}^*(i), i = 1, \dots, N$ 。

(2) 用下式估计 σ^2

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (\hat{\theta}_{ML}^*(i) - \bar{\theta}_{ML}^*)^2 \quad (2.14)$$

$$\bar{\theta}_{ML}^* = \frac{1}{N} \sum_{i=1}^N \hat{\theta}_{ML}^*(i) \quad (2.15)$$

途径 3 用 Taylor 级数展开求 R^* 的一、二阶矩

如前所述, 在 $R(X, F)$ 为对称时, 可记

$$R^* = R(P^*) = R(X^*, \hat{F}_n)$$

于是, 对 $R(P^*)$ 在 P^* 的均值 $P^* = \frac{\sigma}{n}$ 处作 Taylor 展开, 便可得 Bootstrap 分布一、二阶矩

的近似公式。即

$$R(P^*) \cong R\left(\frac{e}{n}\right) + \left(P^* - \frac{e}{n}\right)U + \frac{1}{2}\left(P^* - \frac{e}{n}\right)V\left(P^* - \frac{e}{n}\right)^T \quad (2.16)$$

式中略去了二阶以上项。其中

$$U = \begin{pmatrix} \frac{\partial R^*}{\partial P_1^*} \\ \vdots \\ \frac{\partial R^*}{\partial P_n^*} \end{pmatrix} \Big|_{P^* = \frac{e}{n}} \quad (2.17)$$

$$V = \begin{pmatrix} \frac{\partial^2 R^*}{\partial P_1^* \partial P_1^*} & \cdots & \frac{\partial^2 R^*}{\partial P_1^* \partial P_n^*} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 R^*}{\partial P_n^* \partial P_1^*} & \cdots & \frac{\partial^2 R^*}{\partial P_n^* \partial P_n^*} \end{pmatrix} \Big|_{P^* = \frac{e}{n}} \quad (2.18)$$

注意到 P^* 满足约束 $\sum_{i=1}^n P_i^* = 1$, 则 $eU = 0$, $eV e^T = 0$. 从而由(2.16)有

$$\begin{aligned} E_*[R^*] &\cong R\left(\frac{e}{n}\right) + \frac{1}{2} \text{tr} \left[V \left(\frac{I}{n^2} - \frac{e^T e}{n^3} \right) \right] \\ &= R\left(\frac{e}{n}\right) + \frac{1}{2n} V \end{aligned} \quad (2.19)$$

$$V = \text{tr} \left[V \left(\frac{I}{n} - \frac{e^T e}{n^2} \right) \right] = \text{tr} \frac{V}{n} = \frac{1}{n} \sum_{i=1}^n V_{ii} \quad (2.20)$$

另外, 略去(2.16)中二阶项, 得近似式

$$\begin{aligned} V_{\text{var}_*}[R^*] &= E_*[(R^* - E_*[R^*])^2] = U^T \left(\frac{I}{n^2} - \frac{e^T e}{n^3} \right) U \\ &= \frac{1}{n^2} U^T U = \frac{1}{n^2} \sum_{i=1}^n U_i^2 \end{aligned} \quad (2.21)$$

(2.19)~(2.21)即为所求 R^* 的近似一、二阶矩。

例 3 定义

$$R(X, F) = \theta(\hat{F}_n) - \theta(F) \quad (2.22)$$

其意义在下节将专门讨论。此时

$$R^* = R(X^*, \hat{F}_n) = \theta(\hat{F}_n^*) - \theta(\hat{F}_n) \quad (2.23)$$

其中, \hat{F}_n^* 为 Bootstrap 子样 X^* 的经验分布函数, 即

$$\hat{F}_n^*: \text{在每点 } x_i \text{ 处具有质量 } P_i^*, i=1, \dots, n. \quad (2.24)$$

因此, 当 $P^* = \frac{e}{n}$ 时, 恒有 $\hat{F}_n^* = \hat{F}_n$. 故 $R\left(\frac{e}{n}\right) = 0$.

且

$$E_*[\theta(\hat{F}_n^*) - \theta(\hat{F}_n)] \cong \frac{1}{n}V \quad (2.25)$$

$$V_{a.r.*}[\theta(\hat{F}_n^*) - \theta(\hat{F}_n)] \cong \frac{1}{n^2} \sum_{i=1}^n U_i^2 \quad (2.26)$$

可用上式来近似 $E_r[\theta(\hat{F}) - \theta(F)]$ 和 $V_{a.r.r}[\theta(\hat{F}) - \theta(F)]$ 。

以上讨论了获得 Bootstrap 分布的三种途径。由于途径 1 的适用范围很小, 而途径 3 只能得到 R^* 的近似一、二阶矩, 故实际中最常用的是途径 2。在计算机 (特别是高速计算机) 的支持下, 它可以说是畅通无阻的。这也正是 Bootstrap 方法生命力之所在。

4. 多子样场合的推广

以上就单子样场合讨论了 Bootstrap 方法的基本思想, 对于多子样问题可以类推。以两个子样为例。设 $X = (X_1, \dots, X_m)$ 和 $Y = (Y_1, \dots, Y_n)$ 分别为来自未知总体 F 和 G 的两个相互独立的随机子样, $R((X, Y), (F, G))$ 为预先指定的某一随机变量。在给定观测 $X = x, Y = y$ 之下估计 $R((X, Y), (F, G))$ 的抽样分布特性可采取如下步骤:

(1) 分别构造 X 和 Y 的经验分布函数 \hat{F}_m 和 \hat{G}_n ;

(2) 分别从 \hat{F}_m 和 \hat{G}_n 中抽取 Bootstrap 子样: $X_i^* \sim \hat{F}_m, Y_j^* \sim \hat{G}_n, i = 1, \dots, m, j = 1, \dots, n$;

(3) 计算 $R^* = R((X^*, Y^*), (\hat{F}_m, \hat{G}_n))$ 的分布, 并以此作为 $R((X, Y), (F, G))$ 真实分布的迫近。显然, 上述步骤是单子样时 Bootstrap 方法的自然推广。

三、Bootstrap 方法在统计推断问题中的应用

1. 总体中位数 (median) 的估计

考虑单子样场合。设总体 F 是未知的一维分布, $\theta(F)$ 是 F 的中位数。对于子样 $X = (X_1, \dots, X_n)$, 子样中位数为 $t(X) = X_{(m)}$ 。这里为方便起见设 $n = 2m - 1$ 为奇数, $X_{(1)} \leq \dots \leq X_{(n)}$, 为由子样 X 按大小排列的顺序统计量。在非参数场合通常用 $t(X)$ 作为 $\theta(F)$ 的估计。记

$$R(X, F) = t(X) - \theta(F) = X_{(m)} - \theta(F) \quad (3.1)$$

现要求对 $R(X, F)$ 的分布特性进行统计推断。

仍记来自 \hat{F}_n 的 Bootstrap 子样 $X^* = (X_1^*, \dots, X_n^*)$, N_j^* 为 X_j^* 等于 $x_i (j = 1, \dots, n)$ 的数目, 由前讨论知 $N^* = (N_1^*, \dots, N_n^*)$ 服从多项分布。又注意到对任意正整数 $l, 1 \leq l \leq n$, 事件 $\{X_{(m)}^* > x_{(l)}\}$ 等于事件 $\left\{ \sum_{j=1}^l N_j^* \leq m - 1 \right\}$, 而 $N_{(l)}^* = \sum_{j=1}^l N_j^*$ 服从二项分布 Binomial $\left(n, \frac{l}{n} \right)$ 。于是

$$\begin{aligned} P_*\{X_{(m)}^* > x_{(l)}\} &= P_*\{N_{(l)}^* \leq m - 1\} = P\left\{ \text{Binomial}\left(n, \frac{l}{n}\right) \leq m - 1 \right\} \\ &= \sum_{j=0}^{m-1} \binom{n}{j} \left(\frac{l}{n}\right)^j \left(\frac{n-l}{n}\right)^{n-j} \end{aligned} \quad (3.2)$$

又因为

$$R^* = R(X^*, \hat{F}_n) = X_{(m)}^* - x_{(m)}$$

故有

$$\begin{aligned} P_*\{R^* = x_{(l)} - x_{(m)}\} &= P_*\{X_{(m)}^* = x_{(l)}\} \\ &= P_*\{X_{(m)}^* > x_{(l-1)}\} - P_*\{X_{(m)}^* > x_{(l)}\} \\ &= \sum_{j=0}^{m-1} \left[\binom{n}{j} \left(\frac{l-1}{n}\right)^j \left(\frac{n-l+1}{n}\right)^{n-j} - \binom{n}{j} \left(\frac{l}{n}\right)^j \left(\frac{n-l}{n}\right)^{n-j} \right] \end{aligned} \quad (3.3)$$

可见, 只要给定子样容量 n , 即可得到 $P_*\{R^* = x_{(l)} - x_{(m)}\}$ 的值, 利用它可对 R^* 的 Bootstrap 分布进行估计。例如有

$$E_*[R^*] = \sum_{l=1}^n (x_{(l)} - x_{(m)}) P_*\{R^* = x_{(l)} - x_{(m)}\} \quad (3.4)$$

$$E_*[(R^*)^2] = \sum_{l=1}^n (x_{(l)} - x_{(m)})^2 P_*\{R^* = x_{(l)} - x_{(m)}\} \quad (3.5)$$

这可通过理论计算(途径1)得到。对于 F 具有连续有界的分布密度 f 时, 可以证明估计是渐近一致的。

考虑对 F 的进一步假设, 上述过程可作如下修正:

(1) 已知 F 为对称分布, 可令分布 \hat{F}_{sym} 为

$$\hat{F}_{sym}: \text{在 } x_{(1)}, \dots, x_{(n)} \text{ 和 } 2x_{(m)} - x_{(1)}, \dots, 2x_{(m)} - x_{(n)} \text{ 处具有质量 } \frac{1}{2n-1} \quad (3.6)$$

则 Bootstrap 子样改由 \hat{F}_{sym} 中抽取。若记 $Z_{(1)} \leq Z_{(2)} \leq \dots \leq Z_{(2n-1)}$ 为来自 \hat{F}_{sym} 的顺序子样, 则与(3.3)类似, 对任意正整数 l , $1 \leq l \leq 2n-1$

$$\begin{aligned} P_*\{R^* = Z_{(l)} - x_{(m)}\} &= P\left\{\text{Binomial}\left(n, \frac{l-1}{2n-1}\right) \leq m-1\right\} \\ &\quad - P\left\{\text{Binomial}\left(n, \frac{l}{2n-1}\right) \leq m-1\right\} \end{aligned} \quad (3.7)$$

(2) 若已知 F 为连续分布, 可考虑平滑 Bootstrap。例如, 可把 Bootstrap 子样改为

$$X_i^* = \bar{x} + c(x_{I_i} - \bar{x} + sZ_i), \quad i=1, \dots, n \quad (3.8)$$

其中 \bar{x} , s^2 如前定义; I_i 为从整数集 $\{1, \dots, n\}$ 中独立随机抽取的整数; Z_i 服从某已知的零均值、方差为 σ_z^2 的连续分布 (如 $N(0, \sigma_z^2)$); 而

$$c = (1 + \sigma_z^2)^{-\frac{1}{2}} \quad (3.9)$$

因 $E[x_{I_i}] = \bar{x}$, $V_{or}[x_{I_i}] = s^2$, 易知

$$E_*[X_i^*] = \bar{x} + c(E[x_{I_i}] - \bar{x} + sE[Z_i]) = \bar{x}$$

$$V_{or}[X_i^*] = c^2[V_{or}(x_{I_i}) + s^2V_{or}(Z_i)] = s^2$$

即 X_i^* 与来自 \hat{F}_n 的 Bootstrap 子样具有相同的一、二阶矩。但此时直接的理论计算很困难, 需要通过途径 2。即由(3.8)重复抽样 $x^*(i)$, ($i=1, \dots, N$), 其中 $x^*(i) = (x_1^*(i),$

..., $x_n^*(i)$), 计算得相应的 $R^*(i)$, ($i=1, \dots, n$), 则有

$$E_*[R^*] \cong \frac{1}{N} \sum_{i=1}^N R^*(i) \quad (3.10)$$

$$V_{**}[R^*] \cong \frac{1}{N-1} \sum_{i=1}^n (R^*(i) - E_*[R^*])^2 \quad (3.11)$$

算例[1] 定义

$$R(X, F) = \frac{|t(X) - \theta(F)|}{\sigma(F)}$$

其中 $\sigma(F)$ 为总体 F 的方差。对于 F 为 $N(0, 1)$, $n=13$, 此时 $E_*[R] = 0.95$ 。若设 F 是未知的, 由

$$R^* = R(X^*, \hat{F}_n) = \frac{|t(X^*) - \theta(\hat{F}_n)|}{s}$$

对于 $N=100$, 重复试验 200 次的平均结果见附表 A。可见, 虽 $n=13$ 不大, 结果却比较好。

2. 回归问题

考虑一般的回归模型

$$X_i = g_i(\beta) + \varepsilon_i, \quad i=1, \dots, n \quad (3.12)$$

其中 β 为 $p \times 1$ 未知参数矢量, ε_i ($i=1, \dots, n$) 来自未知分布 F , 且相互独立。设 F 为零均值, 方差 σ_F^2 为未知。所谓回归问题, 即由观测 $x = (x_1, \dots, x_n)$, 求 β 的估计 $\hat{\beta}$, 使 $\sum_{i=1}^n (x_i - g_i(\hat{\beta}))^2 = \min_{\beta} \sum_{i=1}^n (x_i - g_i(\beta))^2$ 。同样, 希望得到有关 $\hat{\beta}$ 抽样分布的一些性质。

运用 Bootstrap 方法解决此问题的步骤如下:

(1) 构造残差 $\hat{\varepsilon}_i = x_i - g_i(\hat{\beta})$ 的经验分布函数

$$\hat{F}_n: \text{在每点 } \hat{\varepsilon}_i - \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i \text{ 处具有质量 } \frac{1}{n}, \quad (i=1, \dots, n) \quad (3.13)$$

(2) 在 \hat{F}_n 中抽取 Bootstrap 子样 $\varepsilon^* = (\varepsilon_1^*, \dots, \varepsilon_n^*)$, 作

$$x_i^* = g_i(\hat{\beta}) + \varepsilon_i^*, \quad i=1, \dots, n \quad (3.14)$$

并以此作为回归模型得到 $\hat{\beta}^*$, 它使

$$\sum_{i=1}^n (x_i^* - g_i(\hat{\beta}^*))^2 = \min_{\beta} \sum_{i=1}^n (x_i^* - g_i(\beta))^2$$

(3) 将上述(2)重复独立进行 N 次, 相应得到 $\hat{\beta}^*(i)$, ($i=1, \dots, N$)。用它们估计 $\hat{\beta}$ 的 Bootstrap 分布, 并以此作为 $\hat{\beta}$ 分布的近似。

考虑线性回归这一特殊情况, 此时有

$$X_i = C_i \beta + \varepsilon_i, \quad i=1, \dots, n$$

写成向量形式有

$$X = C\beta + \varepsilon$$

若记 $C^T C = G$, 显然, $\hat{\beta} = G^{-1} C^T x$, 且 $E[\hat{\beta}] = \beta$, $V_{**}[\hat{\beta}] = \sigma_F^2 G^{-1}$ 。由于 σ_F^2 是未知

的, 其无偏估计为

$$\hat{\sigma}_F^2 = \frac{1}{n-p} \sum_{i=1}^n (x_i - C_i \hat{\beta})^2$$

这些都是众所周知的。

另外同样有

$$E_*[\varepsilon_i^*] = 0, \quad V_{ar}[\varepsilon_i^*] = \frac{1}{n} \sum_{i=1}^n \left(\varepsilon_i - \frac{1}{n} \sum_{i=1}^n \varepsilon_i \right)^2 = \hat{\sigma}^2(\varepsilon_i^*)$$

则

$$\begin{cases} E_*[\hat{\beta}^*] = \hat{\beta} \\ V_{ar_*}[\hat{\beta}^*] = \hat{\sigma}^2(\varepsilon_i^*) G^{-1} \end{cases} \quad (3.15)$$

这里需要指出的是: 在Efron原文[1]构造 \hat{F}_n , \hat{F}_n : 在每点 $\hat{\varepsilon}_i$ 处具有质量 $\frac{1}{n}$, ($i=1, \dots, n$), 时, 认为 \hat{F}_n 为零均值。事实上, 按此定义, $E_*[\varepsilon_i^*] = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i$, 而不是零。这一点应该纠正。

3. 偏倚(Bias)的估计

先引进泛函统计量(Functional Statistic)的概念。设 $\theta(F)$ 是与总体 F 有关的某参数, 它是待估的。称

$$\hat{\theta} = \theta(\hat{F}_n) \quad (3.16)$$

为泛函统计量。它是 $\theta(F)$ 的一个估计。由定义, $\theta(\hat{F}_n)$ 满足对称性条件。常见的泛函统计量有

(1) 对 $\theta(F) = E_F[X]$, $\theta(\hat{F}_n) = \bar{x}$;

(2) 对 $\theta(F) = V_{ar_F}[X]$, $\theta(\hat{F}_n) = s^2$;

(3) 对 $\theta(F) = E_F[Z]/E_F[Y]$, 其中 F 为二维分布, $X = (Y, Z)^T$ 服从 F , $\theta(\hat{F}_n) = \bar{z}/\bar{y}$ 。

若设

$$R(X, F) = \theta(\hat{F}_n) - \theta(F) \quad (3.17)$$

则估计的偏倚为

$$\text{Bias} = E_F[R(X, F)] = E_F[\theta(\hat{F}_n)] - \theta(F) \quad (3.18)$$

下面考虑用Bootstrap方法估计偏倚。因为

$$R^* = R(X^*, \hat{F}_n) = \theta(\hat{F}_n^*) - \theta(\hat{F}_n) \quad (3.19)$$

于是, 作Bootstrap子样 X^* 的经验分布函数

$$\hat{F}_n^*: \text{在每点 } x_i \text{ 处具有质量 } P_i^*, \quad i=1, \dots, n \quad (3.20)$$

从 \hat{F}_n^* 中抽取独立子样得到 $\hat{\theta}^* = \theta(\hat{F}_n^*)$, 重复进行 N 次相应得到 $\hat{\theta}^*(1), \dots, \hat{\theta}^*(N)$ 。于是, 偏倚的估计为

$$\hat{\text{Bias}} = \frac{1}{N} \sum_{j=1}^N \hat{\theta}^*(j) - \hat{\theta} \cong E_*[R^*] \quad (3.21)$$

在例3中曾经简单讨论过此问题。在那里由途径3得到了 $E_*[R^*]$ 的近似公式,忽略二阶以上项后 $E_*[R^*]=0$.这说明 $\hat{B}ias$ 是二阶以上量。

4. 标准偏差(Standard Deviation)估计

所谓标准偏差即误差方差的平方根。该问题已在例2中做了一般性的讨论。这里研究一个具体例子。

设总体 F 为二维分布,它是未知的。现随机子样 $X=(X_1, \dots, X_n)$ 来自 F , $X_i=(U_i, V_i)^T$, U_i 和 V_i 是相关的。要求根据观测 $X=x$ 对 U_i 和 V_i 之间的相关系数 ρ 进行估计。在非参数情况下一般采用Pearrson相关系数

$$\hat{\rho}(x_1, \dots, x_n) = \frac{\sum_{i=1}^n u_i v_i - \sum_{i=1}^n u_i \sum_{j=1}^n v_j / n}{\left\{ \sum_{i=1}^n u_i^2 - \left(\sum_{i=1}^n u_i \right)^2 / n \right\}^{1/2} \left\{ \sum_{i=1}^n v_i^2 - \left(\sum_{i=1}^n v_i \right)^2 / n \right\}^{1/2}} \quad (3.22)$$

这里我们感兴趣的是估计

$$\sigma(\hat{\rho}) = \{E_F[(\rho - \hat{\rho})^2]\}^{1/2} \quad (3.23)$$

对此可运用二维分布时的Bootstrap方法:

(1) 根据观测 $X=x$ 构造经验分布函数:

$$\hat{F}_n: \text{在每点 } x_i = (u_i, v_i)^T \text{ 处具有质量 } \frac{1}{n}, i=1, \dots, n \quad (3.24)$$

(2) 从 \hat{F}_n 中抽取Bootstrap子样 $x^*=(x_1^*, \dots, x_n^*)$, 其中 $x_i^*=(u_i^*, v_i^*)^T$, 并计算 $\hat{\rho}^* = \hat{\rho}(x_1^*, \dots, x_n^*)$;

(3) 重复(2) N 次得到 $\hat{\rho}^*(j)$, $j=1, \dots, N$, 以此构造 $\sigma(\hat{\rho})$ 的估计

$$\hat{\sigma}(\hat{\rho}) = \left\{ \frac{1}{N-1} \sum_{j=1}^N (\hat{\rho}^*(j) - \hat{\rho}^*(\cdot))^2 \right\}^{1/2} \quad (3.25)$$

其中 $\hat{\rho}^*(\cdot) = \frac{1}{N} \sum_{j=1}^N \hat{\rho}^*(j)$ 。当然也可以采用平滑Bootstrap方法, 这里不再详述。

算例[2] 总体 F 为二维正态分布, $X_i=(U_i, V_i)^T$, $E[U_i]=E[V_i]=0$, $V_{\sigma r}[U_i]=V_{\sigma r}[V_i]=1$, $\rho(U_i, V_i)=\frac{1}{2}$ 。取 $n=14$, 则(3.23)真值 $\sigma(\hat{\rho})=0.218$ 。下面对三种形式的Bootstrap过程进行模拟计算, N 为128, 重复试验200次取平均, 结果见附表(B)。

(1) F 未知, 由 \hat{F}_n 中抽取Bootstrap子样;

(2) 设 F 为连续型分布, 考虑均匀平滑Bootstrap, 即从 $\hat{F}_1^{0.5} = \hat{F}_n * (0.5\hat{F}_U)$ 中抽取Bootstrap子样。其中 \hat{F}_U 为一均匀分布, 其协方差阵等于

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

(3) 同(2)考虑正态平滑Bootstrap, 即从 $\hat{F}_2^{0.5} = \hat{F}_n * (0.5\hat{F}_N)$ 抽取Bootstrap子样, 其中 $\hat{F}_N = N_2(\bar{x}, \hat{\Sigma})$ 。此时, 因 F 的真实分布为正态分布, 故采用正态平滑有一定的“利己(self-serving)”性。

由定义, $\hat{F}_1^{0.5}$ 、 $\hat{F}_2^{0.5}$ 和 \hat{F}_n 具有相同相关系数。

5. 判别分析的错分率(Error Rate)估计

设两个独立的随机子样 $X = (X_1, \dots, X_m)$ 和 $Y = (Y_1, \dots, Y_n)$ 分别来自未知的 K 维连续分布 F 和 G 。基于观测 x 和 y , 用线性判别分类的方法可将空间 R^k 分为两个互补的区域 A 和 B , 即

$$\begin{cases} B = \left\{ z: (\bar{y} - \bar{x})^T s^{-1} \left(z - \frac{\bar{x} + \bar{y}}{2} \right) > \log \frac{m}{n} \right\} \\ A = R^k - B \end{cases} \quad (3.26)$$

对于新观测 z , 若 $z \in B$, 则认为 z 来自 G ; 否则认为 z 服从 F 。式中 $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$,

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad s = \frac{1}{m+n} \left[\sum_{i=1}^m (x_i - \bar{x})(x_i - \bar{x})^T + \sum_{i=1}^n (y_i - \bar{y})(y_i - \bar{y})^T \right]。$$

显然, 由于子样的随机性, 上述分类会产生错误。对于分布 F , 定义错分率

$$\text{error}_F = P_F \{ X \in B \} \quad (3.27)$$

其直观的估计是

$$\widehat{\text{error}}_F = \frac{\# \{ x_i \in B \}}{m} \quad (3.28)$$

同样地, 可定义 error_G 和 $\widehat{\text{error}}_G$ 。由问题的对称性, 只需考虑上两式即可。这里我们感兴趣的是

$$R((X, Y), (F, G)) = \widehat{\text{error}}_F - \widehat{\text{error}}_G \quad (3.29)$$

的分布特性。为此运用两子样场合的Bootstrap方法:

(1) 在给定观测 x 和 y 条件下, 构造子样经验分布函数:

$$\hat{F}_m: \text{在每点 } x_i \text{ 处具有质量 } \frac{1}{m}, \quad i = 1, \dots, m \quad (3.30)$$

$$\hat{G}_n: \text{在每点 } y_j \text{ 处具有质量 } \frac{1}{n}, \quad j = 1, \dots, n \quad (3.31)$$

(2) 从 \hat{F}_m 和 \hat{G}_n 中分别抽取独立的Bootstrap子样 x_i^* 和 y_j^* , $i = 1, \dots, m, j = 1, \dots, n$;

(3) 用 \bar{x}^* 、 \bar{y}^* 和 s^* 代替 \bar{x} 、 \bar{y} 和 s , 由(3.26)计算 B^* , 同时计算

$$\begin{aligned} R^* &= R((X^*, Y^*), (\hat{F}_m, \hat{G}_n)) = \widehat{\text{error}}_{\hat{F}_m} - \widehat{\text{error}}_{\hat{F}_n} \\ &= \frac{\# \{ x_i \in B^* \}}{m} - \frac{\# \{ x_i^* \in B^* \}}{m} \end{aligned} \quad (3.32)$$

(4) 重复(2)、(3) N 次获得 $R^*(1), \dots, R^*(N)$ 。用它们可做 R^* 的Bootstrap分布的迫近, 并以此作为 R 真实分布的迫近。例如

$$E_*[R^*] \cong \frac{1}{N} \sum_{j=1}^N R^*(j) \quad (3.33)$$

$$V_{\text{or}_*}[R^*] \approx \frac{1}{N-1} \sum_{j=1}^N (R^*(j) - E_*[R^*])^2 \quad (3.34)$$

算例[1] 设 $F: X \sim N_2\left(\left(-\frac{1}{2}, I\right), I\right)$, $G: Y \sim N_2\left(\left(\frac{1}{2}, I\right), I\right)$ 分别 $m=n=10$ 和 $m=n=20$ 两种情况, 而 $N=100$. (3.29) 式的真实均值和标准偏差由 Monte-Carlo 方法算出. 重复试验 100 次的平均结果如附表 (C) 所示.

6. 比估计 (Ratio Estimation) 问题

设 F 为二维分布, 即有 $X_i = (Y_i, Z_i)^T \sim F$, $i=1, \dots, n$, 且 X_i 相互独立. 比估计即为

$$\theta(F) = \frac{E_F[Y]}{E_F[Z]} \quad (3.35)$$

(设 Y, Z 均大于零). 通常用的统计量为

$$t(X) = \hat{\theta}(F) = \frac{\bar{Y}}{\bar{Z}} = \frac{\sum_{i=1}^n Y_i}{\sum_{i=1}^n Z_i} \quad (3.36)$$

记

$$R(X, F) = \frac{\hat{\theta}(F)}{\theta(F)} \quad (3.37)$$

我们要估计 $R(X, F)$ 的分布特性.

由前面定义可知 $\hat{\theta}(F) = \theta(\hat{F}_n)$ 为泛函统计量. 对此问题的 Bootstrap 解可采用途径 2, 具体作法无特别之处, 故不再详述. 由于 $R(X, F)$ 关于 X 对称, 下面试用途径 3 求出 R^* 的近似一、二阶矩. 注意到

$$R^* = R(X^*, \hat{F}_n) = \frac{t(X^*)}{\theta(\hat{F}_n)} = \frac{\bar{Y}^*}{\bar{Z}^*} \stackrel{\bar{z}}{y} \triangleq R(P^*) \quad (3.38)$$

$$\begin{cases} \bar{Y}^* = \sum_{i=1}^n P_i^* y_i, & \bar{Y}^* |_{P_i^* = \frac{e}{n}} = \bar{y} \\ \bar{Z}^* = \sum_{i=1}^n P_i^* z_i, & \bar{Z}^* |_{P_i^* = \frac{e}{n}} = \bar{z} \end{cases}$$

$$R\left(\frac{e}{n}\right) = \frac{\bar{y}}{\bar{z}} \cdot \frac{\bar{z}}{\bar{y}} = 1$$

故 $\frac{\partial \bar{Y}^*}{\partial P_i^*} = y_i, \quad \frac{\partial \bar{Z}^*}{\partial P_i^*} = z_i$

$$U_i = \frac{\partial R(P^*)}{\partial P_i^*} \Big|_{P_i^* = \frac{e}{n}} = \frac{y_i}{\bar{y}} - \frac{z_i}{\bar{z}}$$

$$V_{i,j} = \frac{\partial^2 R(P^*)}{\partial P_i^* \partial P_j^*} \Big|_{P_i^* = \frac{e}{n}} = \frac{2z_i z_j}{\bar{z}^2} - \frac{y_i z_j + y_j z_i}{\bar{y} \bar{z}}$$

$i, j=1, \dots, n$. 于是

$$\sum_{i=1}^n \frac{z_i}{\bar{z}} = \sum_{i=1}^n \frac{y_i}{\bar{y}} = n$$

$$\begin{aligned} E_*[R^*] &= 1 + \frac{1}{2n} \sum_{i=1}^n V_{ii} = 1 - \frac{1}{n^2} \sum_{i=1}^n \left(\frac{y_i z_i}{\bar{y} \bar{z}} - \frac{z_i^2}{\bar{z} \bar{z}} \right) \\ &= 1 - \frac{1}{n^2} \left\{ \sum_{i=1}^n \left(\frac{y_i}{\bar{y}} - 1 \right) \left(\frac{z_i}{\bar{z}} - 1 \right) - \sum_{i=1}^n \left(\frac{z_i}{\bar{z}} - 1 \right)^2 \right\} \end{aligned} \quad (3.39)$$

$$V_{ar_*}[R^*] = \sum_{i=1}^n \frac{U_i^2}{n^2} = \frac{1}{n^2} \sum_{i=1}^n \left(\frac{y_i}{\bar{y}} - \frac{z_i}{\bar{z}} \right)^2 \quad (3.40)$$

根据以上结果还可以对 $\hat{\theta}(F)$ 作适当的纠偏修正, 得到 $\theta(F)$ 的无偏估计

$$t'(X) = \frac{t(X)}{E_*[R^*]} \quad (3.41)$$

7. 非参数置信区间估计

问题是这样提出的: 设 F 是未知的总体分布, $\theta = \theta(F)$ 为待估的未知分布参数。现以 $\hat{\theta} = \theta(\hat{F}_n)$ 的 Bootstrap 分布为基础, 给出 θ 的置信区间估计的迫近。

下面所用的方法称作“百分值(percentile)”法。记 $\hat{\theta}^* = \theta(\hat{F}_n^*)$, 其 Bootstrap 分布的累积(cumulative)分布函数为

$$\widehat{CDF}(t) = P_{rob_*} \{ \hat{\theta}^* \leq t \} \quad (3.42)$$

对于 Monte—Carlo 模拟近似有

$$\widehat{CDF}(t) \cong \frac{*\{ \hat{\theta}^*(i) \leq t \}}{N} \quad (3.43)$$

给定某一 α , $0 < \alpha < 0.5$, 定义

$$\hat{\theta}_{LOW}(\alpha) = \widehat{CDF}^{-1}(\alpha), \quad \hat{\theta}_{UP}(\alpha) = \widehat{CDF}^{-1}(1 - \alpha) \quad (3.44)$$

简记 $\hat{\theta}_{LOW}(\alpha) = \hat{\theta}_{LOW}$, $\hat{\theta}_{UP}(\alpha) = \hat{\theta}_{UP}$, 有

$$\alpha = \widehat{CDF}(\hat{\theta}_{LOW}), \quad 1 - \alpha = \widehat{CDF}(\hat{\theta}_{UP}) \quad (3.44)'$$

则 $[\hat{\theta}_{LOW}, \hat{\theta}_{UP}]$ 即为在置信水平 $1 - 2\alpha$ 之下 θ 的置信区间估计的迫近。

若 $P_{rob_*} \{ \hat{\theta}^* \leq \hat{\theta} \} = 0.5$, 即 Bootstrap 分布关于 $\hat{\theta}$ 是对称的, 则 θ 的真值近似落在置信区间 $[\hat{\theta}_{LOW}, \hat{\theta}_{UP}]$ 的中间。若对称条件不满足, 则可考虑偏倚校正(bias—corrected)百分值法, 具体作法是: 定义

$$z_0 = \Phi^{-1}(\widehat{CDF}(\hat{\theta})) \quad (3.45)$$

其中 $\Phi(x)$ 为标准正态变量的累积分布函数

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$$

即有

$$\widehat{CDF}(\hat{\theta}) = \Phi(z_0)$$

则偏倚校正百分值法取 z_α , 满足 $\Phi(z_\alpha) = 1 - \alpha$.

$$\hat{\theta}_{LOW} = \widehat{CDF}^{-1}(\Phi(2z_0 - z_\alpha)), \hat{\theta}_{UP} = \widehat{CDF}^{-1}(\Phi(2z_0 + z_\alpha)) \quad (3.46)$$

即

$$\widehat{CDF}(\hat{\theta}_{LOW}) = \Phi(2z_0 - z_\alpha), \widehat{CDF}(\hat{\theta}_{UP}) = \Phi(2z_0 + z_\alpha) \quad (3.46)'$$

则 $[\hat{\theta}_{LOW}, \hat{\theta}_{UP}]$ 为 θ 在置信水平 $1 - 2\alpha$ 之下置信区间的近似表示。特殊地, 当对称条件成立, 显然 $z_0 = 0$ 则 (3.44) 即为 (3.46)。

算例[6] 总体 F 为二维正态分布, 即 $X_i = (U_i, V_i)^T \sim F, i = 1, \dots, n, \rho(U_i, V_i) = 0.5, n = 15$. 若设 F 未在 68% 的置信水平下, $N = 1000$, 重复试验 10 次平均, 关于 P 的置信区间非参数估计的结果见附表(D)。

四、结论与展望

与经典的参数统计理论相比, Bootstrap 方法有两个鲜明的特点:

1. 对于解决某一具体问题, Bootstrap 方法完全从问题的提法本身出发, 仅依赖于给定的观测数据, 而不附加任何其它人为的假设。与此相反, 经典的参数方法则过分依赖于使问题能得到理论解的标准分布模型的小集合及现成的数理统计图表, 这在许多场合难免有些牵强附会。因此 Bootstrap 方法具有更强的现实性。

2. Bootstrap 方法借助于计算机这一有力工具, 使其具有较高的灵活性和实用性。我们已经看到, 无论 $R(X, F)$ 的形式如何复杂, 相应的 Bootstrap 算法都是比较简单的, 也易于在计算机上实现。当然, 从传统的观点看, Bootstrap 方法在计算上存在着巨大的浪费, 它需要将通常的统计计算重复上千次甚至更多。这在计算机没有高度发展的时代是不可想象的, 但在当今做到这点毫不困难。

作为一种新的统计方法, Bootstrap 方法尚在发展中, 许多问题有待于研究解决。这里不妨列举几条:

1. 算法的收敛性问题: Efron 对于一般情况考虑了 Bootstrap 方法的渐近性, 他指出^[1], 当 $R(X, F)$ 的分布弱收敛于某一极限时, $R(X^*, \hat{F}_n)$ 的 Bootstrap 分布以概率 1 弱收敛于同一极限。其它一些学者, 如 Singh^[8]、Bickel 和 Freedman^[9]、Freedman^[10]、Beran^[11] 等人, 都就具体应用问题讨论了 Bootstrap 方法的收敛性, 得到了一些结果。但更一般的收敛性理论还有待于进一步的探讨。

2. 算法的稳健性问题: Beran^[11] 证明在一定条件下抽样分布的 Bootstrap 估计具有 MinMax 特性。另一方面由于 Bootstrap 方法是完全建立在观测子样 $X = x = (x_1, \dots, x_n)$ 的基础之上的。观测数据质量的好坏对 Bootstrap 精度的影响如何? 当子样容量 n 小时 (即小子样场合), Bootstrap 方法的结果又怎样? 这些问题均需要加以回答。

3. 算法的计算量问题: 对于具体问题, 要达到一定的精度要求, 这需要多少计算量? 特别在进行统计模拟计算时, N 应如何选择? 这个问题也值得研究。

其它还有一些问题, 如 Bootstrap 方法与参数统计方法的联系等等。

目前, 在我们的工作中, 参数统计方法仍占着主导地位, 人们习惯于“正态假设”等标准参数模型。Bootstrap 方法的出现为解决实际问题开辟了一条新的途径。我们介绍

Bootstrap 方法的目的是，一是了解在应用统计学领域里的最新发展动态，更重要的是要把这种新方法运用飞行器试验学的具体问题中去，以期提高统计估值的精度和可信度。

最后需要指出的是，Bootstrap方法并不是在一切场合下都适用的。Heran^[11]曾经具体给出了一个反例；Rickel 和Freemdan亦讨论了 Bootstrap 估计失败的问题。这一问题应引起足够的重视。

附 表

(A) 总体中位数的估计

$E_*[R^*]$	无平滑Bootstrap		平滑 Bootstrap, x_i 服从 $[-\frac{d}{2}, \frac{d}{2}]$ 上均匀分布				
	\hat{F}_n	\hat{F}_{sym}	$d=0$	$d=0.25$	$d=0.5$	$d=1$	$d=2$
平均 值	1.01	1.00	1.00	1.01	1.00	0.99	0.93
标准 偏差	0.31	0.33	0.32	0.32	0.32	0.30	0.26

(B) 标准偏差的估计

方 法	类 别	$\sigma_B(\hat{\rho})$	
		平 均	标准 偏差
	无平滑 Bootstrap	0.206	0.066
	正态平滑 Bootstrap	0.200	0.060
	均匀平滑 Bootstrap	0.205	0.061

(C) 错分率的估计

变 量	分 类	$m=n=10$		$m=n=20$	
		均 值	标准偏差	均 值	标准偏差
	R	0.062	0.143	0.028	0.103
	$E_*[R^*]$	0.057	0.026	0.029	0.015
	$\sqrt{V_{\sigma_*}[R^*]}$	0.131	0.016	0.097	0.010

(D) 非参数置信区间估计

方 法	正态理论 (真值)	百分 值 法	偏倚校正百分值法
置信区间	[0.23, 0.33]	[0.27, 0.67]	[0.22, 0.66]

参 考 文 献

- [1] B.Efron, Bootstrap Methods: Another Look at the Jackknife, *Ann. Statist.* 7, 1—36, 1979.
- [2] B.Efron *Computers and the Theory of Statistics: Thinking the Unthinkable*, *SIAM Review* 21, 460—480, 1979.
- [3] B.Efron, Nonparametric Estimates of Standard Error: The Jackknife, the Bootstrap and Other Methods, *Biometrika* 68—3, 589—599, 1981.
- [4] B.Efron, Censored Data and the Bootstrap, *J.American Statistical Association* 76, 312—319, 1981.
- [5] B.Efron, *Computer Intensive Methods in Statistics, Some Recent Advances in Statistics*, Academic Press, 1982.
- [6] B.Efron, *The Jackknife, the Bootstrap and Other Resampling Plans*, CBMS—NSF, SIAM, Philadelphia, 1982.
- [7] B.Efron, A Leisurely Look at the Bootstrap, the Jackknife and the Cross Validation, *Amer. Statist.* 37, 36—48, 1983.
- [8] K.Sigh, On the Asymptotic Theory for the Bootstrap, *Ann. Statist.* 9, 1187—1195, 1981.
- [9] P.J.Bickel & D.A.Freedman, Some Asymptotic Accuracy for the Bootstrap, *Ann. Statist.* 9, 1196—1217, 1981.
- [10] D.A.Freedman, Bootstrapping Regression Models, *Ann. Statist.* 9, 1218—1228, 1981
- [11] R.Beran, Estimated Sampling Distributions: the Bootstrap and Competitors, *Ann. Statist.* 10, 212—225, 1982.
- [12] R.Beran, Jackknife Approximations to Bootstrap Estimates, *Ann. Statist.* 12, 101—118, 1984.

Bootstrap Method and Its Applications

Mao Ning

Abstract

A new statistical inference method, called Bootstrap method, was advanced by Prof. Efron in 1977. The method applies to a variety of estimation problems, so importance has been attached to by more and more statisticians. In the paper the basic idea of the method and its application examples are introduced. And the existing problems and the orientation of the research are discussed as well.